

Density Estimation is a problem of reconstructing the pdf from a set of given data points. Suppose we observe  $x_1, \dots, x_n$  and we want to recover the underlying PDF generating our data set. We assume the data comes from a continuous prob. distribution.

In parametric, we assumed the distribution to be some known distributions and estimate its parameters but here we are not assuming any distribution, rather inventing it.

#### \* Estimating CDF :

Assume that  $X$  follows a continuous distribution with CDF  $F$ . Then, for  $h > 0$   $P(n - \frac{h}{2} < X < n + \frac{h}{2}) = \int_{n-\frac{h}{2}}^{n+\frac{h}{2}} f(y) dy \approx h \cdot f(n)$  (as  $n + \frac{h}{2} - n - \frac{h}{2}$ ) when  $h$  is small.

$$\text{Hence, } \hat{f}(n) = \frac{F(n + \frac{h}{2}) - F(n - \frac{h}{2})}{h} \quad (1).$$

#### \* Empirical CDF :

Let  $x_1, x_2, \dots, x_n$  be iid random variables then the empirical CDF is given by :  $F_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(x_i \leq n) \quad (2).$

$$\hat{f}(n) = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}(n - \frac{h}{2} \leq x_i \leq n + \frac{h}{2}) \quad (\text{EDE}) ;$$

1)  $\hat{F}_n$  is VE of true distribution

2)  $\text{Var}(\hat{F}_n) = \frac{1}{n} F(n)(1 - F(n))$

3)  $\mathbb{P}(\hat{F}_n \leq 1 - \alpha) \approx 1 - \alpha$

3)  $F_n \xrightarrow{?} F(y)$

$$\begin{aligned} E(\hat{F}_n) &= \frac{1}{n} \sum_{i=1}^n E(\mathbb{1}(X_i \leq n)) \\ &= \frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}} \mathbb{1}(X_i \leq n) f(y) dy \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^n f(y) dy \\ &= \frac{n F(y)}{n} \\ &= F(y) \end{aligned}$$

Drawback : A continuous distribution is being estimated by a discrete function.

The estimate requires a large data sample.

\* Def : Kernel Density Estimation (KDE)

A Kernel is called a kernel function if  $K(u) \geq 0$ ,  $\int_{-\infty}^{+\infty} K(u) du = 1$   
 $\& K(-u) = K(u)$  e.g. gaussian kernel.

$$\hat{f}(n) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{n-x_i}{h}\right), \quad h: \text{bandwidth} \quad x_1, \dots, x_n: \text{Sample of } n \text{ observations}$$

However, this is a biased estimator and the variance is complicated.

\* NW Estimator & KDE (Nadaraya & Watson) :

They invented kernel regression :  $\hat{f}(n) = \sum_{i=1}^n y_i w_i$

$$\text{where } w_i = \frac{K\left(\frac{n-n_i}{n}\right)}{\sum_{i=1}^n K\left(\frac{n-n_i}{n}\right)}$$

## \* Kernel Density Estimation & Probabilistic Neural Network

Y Numerical : LR, LASSO, Ridge, Spline, Polynomial, Elastic Net (Regression Models).

Y Categorical : Logistic, Bradley Terry, PNN (Classification models).

Y Count : Poisson (Count Regression)

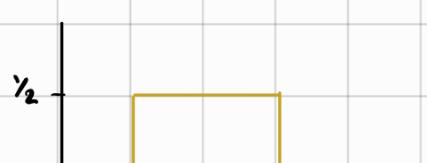
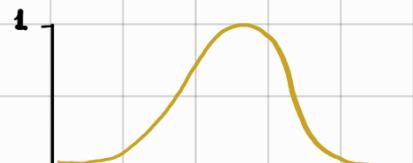
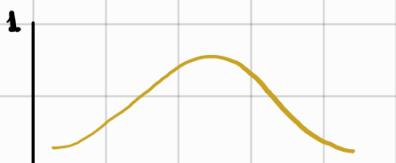
Y Time Series : AR, ARIMA, BSTS, Bass (Forecasting Models)

Y Time to event : Cox PH (Survival Modelling)

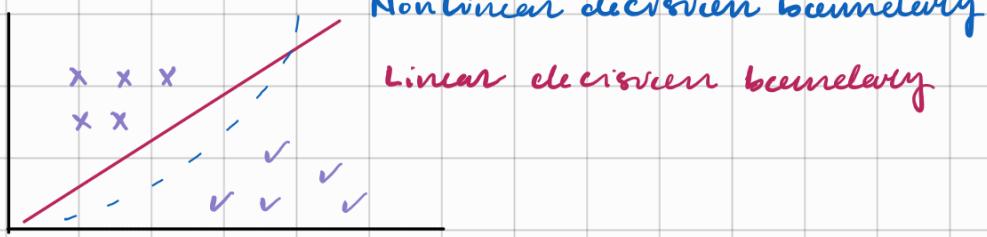
\* Examples : 1) Gaussian Kernel :  $K(u) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{u^2}{2})$

2) Epanechnikov Kernel :  $K(u) = \frac{3}{4} \max\{1-u^2, 0\}$   
(Quadratic)

3) Uniform Kernel :  $K(u) = \begin{cases} 1/2, & |u| < 1 \\ 0, & \text{o.w.} \end{cases}$



- \* PNN :
  - 1) Useful for binary / multi-class classification problem
  - 2) Has nonlinear decision boundary



- 3) Provides probabilistic outputs for classification

#### \* Structure of a probabilistic Neural Network :

- 1) Input Layer : Consists of features of the data
- 2) Pattern Layer : Each neuron in this layer contains one neuron for each training sample. This layer computes a similarity between test input and training samples using kernel.
- 3) Summation Layer : Aggregates the output from pattern layer neurons for each class summing the similarity (kernel values)
- 4) Decision Layer : Calculates prob of test input belonging to each class

Feature 1 ( $x_1$ )	Feature 2 ( $x_2$ )	Output
2.0	2.0	0
2.2	2.2	0 Imbalanced
2.4	2.4	0 data 80% - 20%
2.6	2.6	0
2.8	2.8	0 Imbalance Ratio:
1.8	1.8	0 $IR = \frac{8}{2} = 4$
1.9	2.1	0 $IR = \frac{\# Majority}{\# Minority}$
2.3	2.1	0

4.0

4.0

1

4.1

4.1

1

Test point :  $(x_1, x_2) = (3.5, 3.5)$

$$K(n, x_i) = \exp\left(-\frac{\|n - x_i\|^2}{2\sigma^2}\right)$$

$x_1$	$x_2$	Output	$K(n, x_i)$	Summey	Decision
2.0	2.0	0	0.105	1.98	1.98
2.2	2.2	0	$\exp\left(-\frac{2(3.5-2)^2}{2}\right)$	"	Class 0
2.4	2.4	0	t		
2.6	2.6	0	t		
2.8	2.8	0	t		
1.8	1.8	0	t		
1.9	2.1	0	t		
2.3	2.1	0			
4.0	4.0	1	t	1.47	
4.1	4.1	1			

$$\text{Probability} : P(x \in \{0\}) = \frac{1.98}{1.47 + 1.98} = 0.57$$

Really Close!

$$P(x \in \{1\}) = \frac{1.47}{1.47 + 1.98} = 0.43$$

Fit a Skew Normal to the train

$$\text{Skew normal} : SK(n, x) = K(n, x_i) \Phi\left(\frac{\hat{\alpha} \|n - x_i\|}{\sigma}\right) \quad \hat{\alpha} \in [-6, 6]$$

$$\begin{array}{ll} \text{Skew PNN} : 0 : 0.0347 & P(x \in \{0\}) = 0.16 \\ & P(x \in \{1\}) = 0.84 \\ 1 : 0.185 & \end{array}$$

$\Rightarrow n \in \{1\}$